

# An External-Memory Algorithm for String Graph Construction

Paola Bonizzoni      Gianluca Della Vedova      Yuri Pirola  
 Marco Previtali      Raffaella Rizzi

DISCo, Univ. Milano-Bicocca, Milan, Italy

Some recent results [13, 25] have introduced external-memory algorithms to compute self-indexes of a set of strings, mainly via computing the Burrows-Wheeler Transform (BWT) of the input strings. The motivations for those results stem from Bioinformatics, where a large number of short strings (called reads) are routinely produced and analyzed. In that field, a fundamental problem is to assemble a genome from a large set of much shorter samples extracted from the unknown genome. The approaches that are currently used to tackle this problem are memory-intensive. This fact does not bode well with the ongoing increase in the availability of genomic data. A data structure that is used in genome assembly is the string graph, where vertices correspond to samples and arcs represent two overlapping samples. In this paper we address an open problem of [27]: to design an external-memory algorithm to compute the string graph.

## 1 Introduction

Several fields are witnessing a huge increase in the amount of available data, such as real-time network data, Web usage logs, telephone call records, financial transactions, and biological data [7, 11]. There are three main algorithmic solutions to cope with that amount of data: (1) data streaming, where only one pass is made over the data and the working memory is small compared to the input data [15, 18], (2) parallel algorithms, where input data is split among several processors [30], and (3) external memory algorithms [2, 32] where only part of the data is kept in main memory and most of the data are on disk.

The latter subfield has blossomed with the seminal paper by Vitter and Shriver [31], introducing the *parallel disk model*, where the performance is measured as the number of I/O operations and the amount of disk space used.

The above discussion is especially relevant in Bioinformatics, where we are currently witnessing a tremendous increase in the data available, mainly thanks to the rise of

different Next-Generation Sequencing (NGS) technologies [21]. De novo sequence assembly is still one of the most fundamental problems and is currently receiving a lot of attention, just as it used to be twenty years ago [3,4]. The assembly problem asks for a superstring  $G$  (the unknown genome) of the set  $R$  of input strings (sampled from the unknown genome). The concatenation of all input strings is a feasible solution of the problem, but it is clearly a solution void of any biological significance: for this reason a suitable optimization criterion must be introduced. The simplest criterion is to find a shortest superstring of  $R$  [9], but that model considers neither that the input strings are sampled uniformly from the genome  $G$ , nor that the samples may contain some errors (that is  $G$  is a superstring only in an approximate sense). Moreover, data obtained with different technologies or different instruments have different characteristics, such as the length of the samples and the error distribution, making difficult to describe a unified computational problem that actually represents the real-world genome assembly problem.

For all those reasons, the successful assemblers incorporate a number of ideas and heuristics originating from the biological characteristics of the input data and of the expected output. Interestingly, almost all the assemblers are based on some notion of graph to construct a draft assembly. Most of the available assemblers [5,24,29] are built upon the notion of de Bruijn graphs, where the vertices are all distinct  $k$ -mers (that is the  $k$ -long substrings appearing in at least an input string). If we want to analyze datasets coming from different technologies, hence with important variability in read lengths, an approach based on same-length strings is likely to be limiting. Moreover, one of the main hurdles to overcome is the main memory space that is used by those assemblers. For instance, a standard representation of the de Bruijn graph for the human genome when  $k = 27$  requires 15GB (and is unfeasible in metagenomics). To reduce the memory usage, a probabilistic version of de Bruijn graphs, based on the notion of Bloom filter, has been introduced [12] and uses less than 4GB of memory to store the de Bruijn graph for the human genome when  $k = 27$ .

The amount of data necessary to assemble a genome emphasizes the need for algorithmic solutions that are time and space efficient. An important challenge is to reduce main memory usage while keeping a reasonable time efficiency. For this reason, some alternative approaches have been developed recently, mostly based on the idea of *string graph*, initially proposed by Myers [22] before the advent of NGS technologies, and further developed [27,28] to incorporate some advances in text indexing, such as the FM-index [17]. These methods build a graph whose vertices are the input reads and a visit of the paths of the graph allows to reconstruct the genome.

A practical advantage of string graphs over de Bruijn graphs is that reads are usually much longer than  $k$ , therefore string graphs can immediately disambiguate some repeats that de Bruijn methods might resolve only at later stages. On the other hand, string graphs are more computationally intensive to compute [28]. For this reason we have studied the problem of computing the string graph on a set  $R$  of input strings, with the goal of developing an external-memory algorithm that uses only a limited amount of main memory, while minimizing disk accesses.

Our work has been partially inspired by SGA [27], the most used string graph assembler. SGA, just as several other bioinformatics programs, is based on the notions of BWT [10]

and of FM-index constructed from the set  $R$  of reads. In fact, an important distinguishing feature of SGA is its use of the FM-index to compute the arcs of the string graph. Still, the memory usage of SGA is considerable, since the experimental analysis in [28] has proved that SGA can successfully assemble the human genome from a set of  $\approx 1$  billion 101bp reads, but uses more than 50GB of RAM to complete the task. The space improvement achieved in the latest SGA implementation [28] required to apply a distributed construction algorithm of the FM-index for the collection of reads and specific optimization strategies to avoid keeping the whole indexing of reads in main memory. Indeed, the authors of SGA [27] estimated 400GB of main memory to build the FM-index for a collection of reads at 30x coverage over the human genome. Since the approach used in SGA [27] requires to keep in main memory the entire BWT and the FM-index of all input data, an open problem of [27] is to reduce the space requirements by developing an external memory algorithm to compute the string graph. In this paper we are going to address this open problem.

Another fundamental inspiration is the sequence of papers [6, 14, 25] that have culminated in BCReXt [6], a lightweight (*i.e.*, external-memory) algorithm to compute the BWT (as well as a number of other data structures) of a set of strings. In fact, our algorithm receives as input all data structures computed by BCReXt and the set  $R$  of input reads, computing the string graph of  $R$ .

## 2 Definitions

We briefly recall some standard definitions that will be used in the following. Let  $\Sigma$  be an ordered finite alphabet and let  $S$  be a string over  $\Sigma$ . We denote by  $S[i]$  the  $i$ -th symbol of  $S$ , by  $l = |S|$  the length of  $S$ , and by  $S[i : j]$  the substring  $S[i]S[i + 1] \cdots S[j]$  of  $S$ . The *reverse* of  $S$  is the string  $S^r = S[l]S[l - 1] \cdots S[1]$ . The *suffix* and *prefix* of  $S$  of length  $k$  are the substrings  $S[l - k + 1 : l]$  (denoted by  $S[l - k + 1 : ]$ ) and  $S[1 : k]$  (denoted by  $S[: k]$ ) respectively. The  $k$ -suffix of  $S$  is the  $k$ -long suffix of  $S$ .

Given two strings  $(S_i, S_j)$ , we say that  $S_i$  *overlaps*  $S_j$  iff a nonempty suffix  $Z$  of  $S_i$  is also a prefix of  $S_j$ , that is  $S_i = XZ$  and  $S_j = ZY$ . In that case we say that that  $Z$  is the *overlap* of  $S_i$  and  $S_j$ , denoted as  $ov_{i,j}$ , that  $Y$  is the *right extension* of  $S_i$  with  $S_j$ , denoted as  $rx_{i,j}$ , and  $X$  is the *left extension* of  $S_j$  with  $S_i$ , denoted as  $lx_{i,j}$ .

In the following of the paper we will consider a collection  $R = \{r_1, \dots, r_m\}$  of  $m$  strings (also called *reads*, borrowing the term from the bioinformatics literature) over  $\Sigma$ . As usual, we append a sentinel symbol  $\$ \notin \Sigma$  to the end of each string ( $\$$  lexicographically precedes all symbols in  $\Sigma$ ) and we denote by  $\Sigma^\$$  the extended alphabet  $\Sigma \cup \{\$\}$ . We assume that the sentinel symbol  $\$$  is not taken into account when computing overlaps between two strings or when considering the length of a string. Moreover, two sentinel symbols do not match when compared with each other. This fact implies that two strings consisting only of a sentinel symbol have a longest common prefix that has length zero. A technical difficulty that we will overcome is to identify each read even using a single sentinel. Using a distinct sentinel for each read would eliminate the problem, but it would make the alphabet size too large. We denote by  $n$  the total number of characters

in the input strings, and by  $l$  the maximum length of a string, that is  $n = \sum_{i=1}^m |r_i|$  and  $l = \max_{i=1}^m \{|r_i|\}$ .

**Definition 1.** The *Generalized Suffix Array (GSA)* [26] of  $R$  is the array  $SA$  where each element  $SA[i]$  is equal to  $(k, j)$  if and only if the  $k$ -suffix of string  $r_j$  is the  $i$ -th smallest element in the lexicographic ordered set of all suffixes of the strings in  $R$ . The *Longest Common Prefix (LCP)* array of  $R$ , is the  $n$ -long array  $L$  such that  $L[i]$  is equal to the length of the longest prefix shared by the the  $k_i$ -suffix of  $r_{j_i}$  and the  $k_{i-1}$ -suffix of  $r_{j_{i-1}}$ , where  $SA[i] = (k_i, j_i)$  and  $SA[i-1] = (k_{i-1}, j_{i-1})$ . Conventionally,  $L[1] = -1$ . The *Burrows-Wheeler Transform (BWT)* of  $R$  is the sequence  $B$  such that  $B[i] = r_j[k-1]$ , if  $SA[i] = (k, j)$  and  $k > 1$ , or  $B[i] = \$$ , otherwise. Informally,  $B[i]$  is the symbol that precedes the  $k$ -suffix of string  $r_j$  where such suffix is the  $i$ -th smallest suffix in the ordering given by  $SA$ .

$i$	$LS[i]$	$SA[i]$	$L[i]$	$B[i]$
1	\$	(0, 1)	-	E
2	\$	(0, 3)	0	T
3	\$	(0, 2)	0	N
4	APPLE\$	(5, 1)	0	\$
5	APRICOT\$	(7, 3)	2	\$
6	COT\$	(3, 3)	0	I
7	E\$	(1, 1)	0	L
8	EMON\$	(4, 2)	1	L
9	ICOT\$	(4, 3)	0	R
10	LE\$	(2, 1)	0	P
11	LEMON\$	(5, 2)	2	\$
12	MON\$	(3, 2)	0	E
13	N\$	(1, 2)	0	O
14	ON\$	(2, 2)	0	M
15	OT\$	(2, 3)	1	C
16	PLE\$	(3, 1)	0	P
17	PPLE\$	(4, 1)	1	A
18	PRICOT\$	(6, 3)	1	A
19	RICOT\$	(5, 3)	0	P
20	T\$	(1, 3)	0	O

Table 1: GSA, LCP, BWT on the reads APPLE, LEMON, APRICOT

The  $i$ -th smallest (in lexicographic order) suffix is denoted  $LS[i]$ , that is if  $SA[i] = (k, j)$  then  $LS[i] = r_j[|r_j| - k + 1 :]$ . In the paper, and especially in the statements, we assume that  $R$  is a set of reads,  $SA$  is the generalized suffix array of  $R$ ,  $L$  is the LCP array of  $R$ , and  $B$  is the Burrows-Wheeler Transform of  $R$ .

Given a string  $Q$  and a collection  $R$ , notice that all suffixes of  $R$  whose prefix is  $Q$  appear consecutively in  $LS$ . We call  $Q$ -interval [6] on  $R$  (or simply  $Q$ -interval, if the set

$R$  is clear from the context) the maximal interval  $[b, e)$  such that  $Q$  is a prefix of  $LS[i]$  for each  $i$ ,  $b \leq i < e$  (we denote the  $Q$ -interval by  $q(Q)$ ). Sometimes we will need to refer to the  $Q$ -interval on the set  $R^r$ : in that case the  $Q$ -interval is denoted by  $q^r(Q)$ . We define the *length* and *width* of the  $Q$ -interval  $[b, e)$  on  $R$  as  $|Q|$  and the difference  $e - b$ , respectively. Notice that the width of the  $Q$ -interval is equal to the number of occurrences of  $Q$  as a substring of some string  $r \in R$ . For instance, on the example in Table 1 the LE-interval is  $[10, 12)$ . Whenever the string  $Q$  is not specified, we will use the term *string-interval*. Since the BWT, the LS, the GSA, and the LCP arrays are all closely related, a string interval can be viewed as an interval on any of those arrays.

To extend the previous definition of string interval to consider a string  $Q$  that is a string over the alphabet  $\Sigma^\$$ , we have some technical details to fix, related to the fact that a suffix can contain a sentinel  $\$$  only as the last character. For example, we have to establish what is the  $\$A$ -interval in Table 1, even though no suffix has  $\$A$  as a prefix. To overcome this hurdle, to each suffix  $S[i:]$  of the string  $S$  we associate  $S[i:]S[i-1]$ , which is a rotation of  $S$ , and let  $LS'$  be array of the sorted rotations. Given a string  $Q$  over  $\Sigma^\$$ , we define the  $Q$ -interval as the maximal interval  $[b, e)$  such that  $Q$  is a prefix of the  $i$ -th rotation (in lexicographic order) for each  $i$ ,  $b \leq i < e$ .

Let  $S$  be a string over  $\Sigma$ . Then the  $S\$$ -interval  $q(S\$)$  contains exactly one suffix extracted from each read with suffix  $S$ . Moreover, the  $\$S$ -interval  $q(\$S)$  contains exactly one suffix extracted from each read with prefix  $S$ . For this reason, we will say that  $q(S\$)$  identifies the set  $R^s(S)$  of the reads with suffix  $S$ , and  $q(\$S)$  identifies the set  $R^p(S)$  of the reads with prefix  $S$ .

Let  $B^r$  be the BWT of the set  $R^r = \{r^r \mid r \in R\}$ , let  $[b, e)$  be the  $Q$ -interval on  $R$  for some string  $Q$ , and let  $[b^r, e^r)$  be the  $Q^r$ -interval on  $R^r$ . Then,  $[b, e)$  and  $[b^r, e^r)$  are called *linked* string-intervals. The linking relation is a 1-to-1 correspondence and two linked intervals have same width and length, hence  $e - b = e^r - b^r$ . Given two strings  $Q_1$  and  $Q_2$ , the  $Q_1$ -interval and the  $Q_2$ -interval on  $R$  are either contained one in the other (possibly are the same) or disjoint. There are some interesting relations between string-intervals and the LCP array.

The interval  $[i : j]$  of the LCP array is called an lcp-interval of value  $k$  (shortly  $k$ -interval) if  $L[i], L[j+1] < k$ , while  $L[h] \geq k$  for each  $h$  with  $i < h \leq j$  and there exists  $L[h] = k$  with  $i < h \leq j$  [1]. An immediate consequence is the following proposition.

**Proposition 1.** *Let  $R$  be a set of reads, let  $L$  be the LCP array of  $R$ , let  $S$  be a string and let  $[b, e)$  be the  $S$ -interval. Then  $L[h] \geq |S|$  for each  $h$  with  $b < h \leq e - 1$ . Moreover, if  $S$  is the longest string whose  $S$ -interval is  $[b, e)$ , then  $[b : e - 1]$  is a  $|S|$ -interval.*

Proposition 1 relates the notion of string-intervals with that of lcp-intervals. It is immediate to associate to each  $k$ -interval  $[b, e)$  the string  $S$  consisting of the common  $k$ -long prefix of all suffixes in  $LS[i]$  with  $b \leq i < e$ . Such string  $S$  is called the *representative* of the  $k$ -interval. Moreover, given a  $k$ -interval  $[b, e)$ , we will say that  $b$  is its *opening position* and that  $e$  is its *closing position*.

**Proposition 2.** *Let  $S_1, S_2$  be two strings such that the  $S_2$ -interval  $[b_2, e_2)$  is nonempty, and let  $[b_1, e_1)$  be the  $S_1$ -interval. Then  $S_1$  is a proper prefix of  $S_2$  if and only if  $[b_1, e_1)$  contains  $[b_2, e_2)$  and  $|S_1| < |S_2|$ .*

*Proof.* The only if direction is immediate, therefore we only consider the case when  $[b_1, e_1)$  contains  $[b_2, e_2)$  and  $|S_1| < |S_2|$ .

If the containment is proper, the proof is again immediate from the definition of string-interval. Therefore assume that  $b_1 = b_2$  and  $e_1 = e_2$ . Since  $[b_1, e_1)$  is nonempty, the representative  $S$  of  $[b_1, e_1)$  is not the empty string. Moreover both  $S_1$  and  $S_2$  are prefixes of  $S$ , since  $S$  is the longest common prefix of  $LS[b, e]$ . Since  $|S_1| < |S_2|$ ,  $S_1$  is a proper prefix of  $S_2$ .  $\square$

Notice that Proposition 2 is restricted to nonempty  $S_2$ -intervals, since the  $Q$ -interval is empty for each string  $Q$  that is not a substring of a read in  $R$ . Therefore relaxing that condition would falsify Proposition 2. On the other hand, we do not need to impose the  $S_1$ -interval to be nonempty, since it is an immediate consequence of the assumption that  $S_1$  is a prefix of  $S_2$ .

Given a  $Q$ -interval and a symbol  $\sigma^\$ \in \Sigma$ , the *backward  $\sigma$ -extension* of the  $Q$ -interval is the  $\sigma Q$ -interval (that is, the interval on the GSA of the suffixes sharing the common prefix  $\sigma Q$ ). We say that a  $Q$ -interval has a *nonempty* (*empty*, respectively) backward  $\sigma$ -extension if the resulting interval has width greater than 0 (equal to 0, respectively). Conversely, the *forward  $\sigma$ -extension* of a  $Q$ -interval is the  $Q\sigma$ -interval. We recall that the FM-index [17] is essentially made of the two arrays  $C$  and  $Occ$ , where  $C(\sigma)$ , with  $\sigma \in \Sigma$ , is the number of occurrences in  $B$  of symbols that are alphabetically smaller than  $\sigma$ , while  $Occ(\sigma, i)$  is the number of occurrences of  $\sigma$  in the prefix  $B[: i - 1]$  (hence  $Occ(\cdot, 1) = 0$ ). It is immediate to obtain from  $C$  a function  $C^{-1}(i)$  that returns the first character of  $LS[i]$ ; this fact allows to represent a character with an integer. We can now state a fundamental characterization of extensions of a string-interval [17, 20]. This characterization allows to compute efficiently all extensions, via  $C$  and  $Occ$ .

**Proposition 3.** *Let  $S$  be a string, let  $q(S) = [b, e)$  be the  $S$ -interval and let  $\sigma$  be a character. Then the backward  $\sigma$ -extension of  $[b, e)$  is  $q(\sigma S) = [C(\sigma) + Occ(\sigma, b) + 1, C(\sigma) + Occ(\sigma, e))$ , and the forward  $\sigma$ -extension of  $[b, e)$  is  $q(S\sigma) = [b + \sum_{c < \sigma} (Occ(\sigma, e) - Occ(\sigma, b)), b + \sum_{c \leq \sigma} (Occ(\sigma, e) - Occ(\sigma, b))$ .*

Proposition 3 presents a technical problem when  $\sigma$  is the sentinel  $\$$ . More precisely, since all reads share the same sentinel  $\$$ , we might not have a correspondence between suffixes in the  $q(\$S)$  and reads with prefix  $S$ . More precisely, if  $q(\$S) = [b, e)$ ,  $b \leq i < e$ , and  $SA[i] = (k, j)$ , we do not know whether  $r_j$  has prefix  $S$  (which is needed to preserve the correspondence between  $q(\$S)$  and  $R^p(S)$ ). For this reason, we sort the suffixes that are equal to the sentinel  $\$$  (corresponding to the positions  $i$  such that  $SA[i] = (0, \cdot)$ ) according to the lexicographic order of the reads. In other words, we enforce that, for each  $i_1, i_2$  where  $SA[i_1] = (0, j_1)$ ,  $SA[i_2] = (0, j_2)$ , if  $i_1 < i_2$  then  $r_{j_1}$  lexicographically precedes  $r_{j_2}$ . For that purpose, it suffices a coordinated scan of the GSA and of the BWT, exploiting the fact that  $B[i] = \$$  and  $SA[i] = (k, j)$  iff the read  $r_j$  is  $k$  long.

**Definition 2** (Overlap graph). Given a set  $R$  of reads, its *overlap graph* [22] is the directed graph  $G_O = (R, A)$  whose vertices are the reads in  $R$ , and where two reads  $r_i, r_j$  form the arc  $(r_i, r_j)$  if they have a nonempty overlap. Moreover, each arc  $(r_i, r_j)$  of  $G_O$  is labeled by the left extension  $lx_{i,j}$  of  $r_i$  with  $r_j$ .

The main use of string graph is to compute the *assembly* of each path, corresponding to the sequence that can be read by traversing the reads corresponding to vertices of the path and overlapping those reads. More formally, given a path  $r_{i_1}, r_{i_2}, \dots, r_{i_k}$  of  $G_O$ , its assembly is the string  $lx_{i_1, i_2} lx_{i_2, i_3} \dots lx_{i_{k-1}, i_k} r_{i_k}$ .

The original definition of overlap graph [22] differs from ours since the label of the arc  $(r_i, r_j)$  consists of the right extension  $rx_{i,j}$  as well as  $lx_{i,j}$ . Accordingly, also their definition of assembly uses the right extensions instead of the left extensions. The following lemma establishes the equivalence of those two definitions in terms of *assembly* of a path.

**Lemma 4.** *Let  $G_O$  be the overlap graph for  $R$  and let  $r_{i_1}, r_{i_2}, \dots, r_{i_k}$  be a path of  $G_O$ . Then,  $lx_{i_1, i_2} lx_{i_2, i_3} \dots lx_{i_{k-1}, i_k} r_{i_k} = r_{i_k} rx_{i_1, i_2} rx_{i_2, i_3} \dots rx_{i_{k-1}, i_k}$ .*

*Proof.* We will prove the lemma by induction on  $k$ . Let  $(r_h, r_j)$  be an arc of  $G_O$ . Notice that the path  $r_h r_j$  represents  $lx_{h,j} ov_{h,j} rx_{h,j}$ . Since  $r_h = lx_{h,j} ov_{h,j}$  and  $r_j = ov_{h,j} rx_{h,j}$ , the case  $k = 2$  is immediate.

Assume now that the lemma holds for paths of length smaller than  $k$  and consider the path  $(r_{i_1}, \dots, r_{i_k})$ . The same argument used for  $k = 2$  shows that  $lx_{i_1, i_2} lx_{i_2, i_3} \dots lx_{i_{k-1}, i_k} r_{i_k} = lx_{i_1, i_2} lx_{i_2, i_3} \dots lx_{i_{k-2}, i_{k-1}} r_{i_{k-1}} rx_{k-1, k}$ . By inductive hypothesis  $lx_{i_1, i_2} lx_{i_2, i_3} \dots lx_{i_{k-2}, i_{k-1}} r_{i_{k-1}} rx_{k-1, k} = r_{i_1} rx_{i_1, i_2} rx_{i_2, i_3} \dots rx_{i_{k-2}, i_{k-1}} rx_{i_{k-1}, i_k}$ , completing the proof.  $\square$

This definition models the actual use of string graphs to reconstruct a genome [22]. If we have perfect data and no relevant repetitions, the overlap graph is a directed acyclic graph (DAG) with a unique topological sort, which in turn reveals a peculiar structure; the graph is made of tournaments [16]. More formally, let  $\langle r_1, \dots, r_n \rangle$  be the topological order of  $G_O$ . If  $(r_i, r_j)$  is an arc of  $G_O$  then also all  $(r_h, r_k)$  with  $i < h < k < j$  are arcs of  $G_O$ . Notice that in this case, all paths from  $r_i$  to  $r_j$  have the same assembly.

Less than ideal conditions might violate the previous property. In fact insufficient coverage (where we do not have reads extracted from some parts of the original genome) or sequencing errors (where the read is not a substring of the genome) might result in a disconnected graph, while spurious overlaps or long repetitions might result in a graph that is not a DAG. Nonetheless, the ideal case points out that we can have (and we actually have in practice) multiple paths with the same assembly.

This suggests that it is possible (and auspicious) to remove some arcs of the graph without modifying the set of distinct assemblies. An arc  $(r_i, r_j)$  of  $G_O$  is called *reducible* [22] if there exists another path from  $r_i$  to  $r_j$  with the same assembly (*i.e.*, the string  $lx_{i,j} r_j$ ). After removing all reducible arcs we obtain the *string graph* [22].

In this paper we are going to develop two external-memory algorithms, the first to compute the overlap graph associated to a set of reads, and the second to reduce an overlap graph into a string graph.

For simplicity, and to emphasize that our algorithms are suited also for an in-memory implementation, we use lists as main data structures. An actual external-memory implementation will replace such lists with files that can be accessed only sequentially. We will use an array-like notation to denote each element, but accessing those elements

sequentially. Moreover, we will assume that the set of reads  $R$  has been processed with the BCRext algorithm [6] to compute the BWT, GSA and LCP of  $R$ .

### 3 Computing the overlap graph

Our algorithm for computing the overlap graph is composed of two main parts: (i) computing the unlabeled overlap graph, (ii) labeling the arcs. Notice that, given a string  $S$ , the cartesian product  $R^s(S) \times R^p(S)$ , where  $R^p(S)$ ,  $R^s(S)$  are respectively the set of reads in  $R$  whose prefix and suffix (respectively) is  $S$ , consists of the arcs whose overlap is  $S$ . Observe that the pair  $(q(S\$), q(\$S))$  of string-intervals represents the set of arcs whose overlap is  $S$ , since  $q(S\$)$  and  $q(\$S)$  represent the sets  $R^s(S)$  and  $R^p(S)$ , respectively. Characterizing also the arc labels is more complicated, as pointed out by Definition 3.

A consequence of Lemma 4 is that we can label each arc with its left extension. Indeed, given a read  $r_i = PS$  we use the  $P$ -interval to label the arcs  $(r_i, r_j)$  (outgoing from  $r_i$ ) with overlap  $ov_{i,j} = S$  and left extension  $lx_{i,j} = P$ . Anyway, to compute the  $P$ -interval we will need also the  $PS\$$ -interval. Moreover, our procedure that reduces an overlap graph is based on a property relating the reverse  $P^r$  of the left extension  $P$ ; for this reason we need to encode  $P^r$  as well as  $P$ .

**Definition 3.** Let  $P$  and  $S$  be two strings. Then, the tuple  $(q(PS\$), q(\$S), q(P), q^r(P^r), |P|, |S|)$  is the  $(P, S)$ -*encoding* (or simply *encoding*) of all arcs with left extension  $P$  and overlap  $S$ . Moreover, the  $(P, S)$ -encoding is *terminal* if the  $PS\$$ -interval has a nonempty backward  $\$$ -extension, and it is *basic* if  $P$  is the empty string.

Notice that a basic  $(\epsilon, S)$ -encoding is equal to  $(q(S\$), q(\$S), q(\epsilon), q(\epsilon), 0, |S|)$ , where the interval  $q(\epsilon)$  is  $[1, n + m + 1)$ , where  $n + m$  is the overall number of characters in the input reads, included the sentinels. Moreover, the differences between basic and non-basic encodings consist of the information on  $P$ , that is the arc label. In other words, the basic encodings already represent the topology of the overlap graph. For this reason, the first part of our algorithm will be to compute all basic encodings. Moreover, we want to read sequentially three lists— $\mathcal{B}$ ,  $\mathcal{SA}$  and  $\mathcal{L}$ —that have been previously computed via BCRext [6] containing the BWT  $B$ , the GSA  $SA$  and the LCP array  $L$ , respectively. Another goal of our approach is to minimize the number of passes over those lists, as a simpler adaptation of the algorithm of [27] would require a number of passes equal to the sum of the lengths of the input reads in the worst case, which would clearly be inefficient.

**Definition 4.** Let  $S$  be a proper substring of some read of  $R$ . If both the  $S\$$ -interval and the  $\$S$ -interval are nonempty, then  $S$  is called a *seed*.

Finding all basic encodings is mostly equivalent to finding all seeds of  $R$ . Now we will prove that a position  $p$  can be the opening position of several lcp-intervals, but it can be the opening position of only one seed.

**Lemma 5.** Let  $S$  be a  $k$ -long seed and let  $[b, e)$  be the corresponding  $k$ -interval. Then  $k = \mathcal{L}[b + 1]$  and  $k > \mathcal{L}[b]$ .



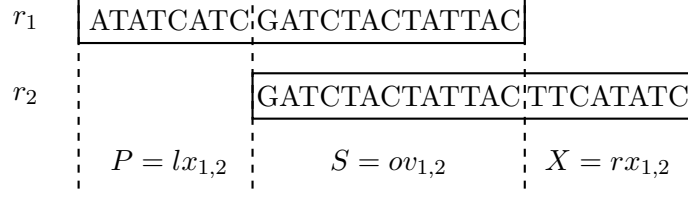


Figure 1: Example of arc of the overlap graph. The read  $r_1$  is equal to ATATCATCGATCTACTATTAC, while the read  $r_2$  is equal to GATCTACTATTACTTTCATATC.

*Proof.* Since  $S$  is a seed,  $S$  has a forward  $\$$ -extension. Moreover  $\$$  is the smallest character, hence  $S\$$  is the first suffix of the  $S$ -interval, that is the  $LS[b] = S\$$ . Since  $S$  is a seed,  $S$  is not an entire read. Hence the  $\$S$ -interval and the  $S\$$ -interval are both nonempty and are disjoint. Consequently,  $S$  is a prefix of  $LS[b+1]$ , that is  $LS[b+1] = S\alpha$  for some string  $\alpha$  on the alphabet  $\Sigma \cup \{\$\}$ . By definition,  $\mathcal{L}[b+1]$  is the length of the longest common prefix of  $S\$$  and  $S\alpha$ , that is  $\mathcal{L}[b+1] = |S|$ .

If  $b = 1$ , by definition  $\mathcal{L}[b] = -1$ , hence  $\mathcal{L}[b] < |S|$ . If  $b > 1$ ,  $S$  is not a prefix of  $LS[b]$ , hence  $\mathcal{L}[b] < |S|$ .  $\square$

**Corollary 6.** *Let  $[b, e)$  be an interval. Then  $b$  is the opening position of at most one seed.*

In the presentation of our algorithm, we need a simple procedure called **Merge**. This procedure operates on lists of encodings  $([b, e), \cdot, \cdot, \cdot, \cdot)$  that are sorted by increasing  $b$  (records with the same  $b$  are sorted by decreasing  $e$ ). We do not actually need to write a new list, as **Merge** consists of choosing the list from which to read the next record.

To compute all basic encodings, the procedure **BuildBasicArcIntervals** (Algorithm 1) reads sequentially the lists  $\mathcal{B}$ ,  $\mathcal{L}$  and  $\mathcal{SA}$  while keeping in  $\#_\$$  the number of sentinels in the portion of BWT before the current position. Moreover we maintain a stack  $Z$  that is used to store the relevant  $k$ -intervals whose opening position has been read, but whose closing position has not.

When the current position is  $p$ , the only interesting cases are if  $p - 1$  is an opening position or  $p$  is a closing position. In the first case (see lines 14–23), Lemma 5 and Corollary 6 show that only the  $\mathcal{L}[p]$ -long interval whose opening position is  $p - 1$  might have a seed as representative. Let  $S$  be the representative of such  $\mathcal{L}[p]$ -long interval. First, at lines 19–23, we compute the  $S\$$ -interval  $[p - 1, e_1)$  (since such string interval is an initial portion of the interval, we only need to read some records from the input lists). If the  $S\$$ -interval is nonempty, then it is pushed onto  $Z$  together with the current value of  $\#_\$$  and the length of  $S$ . Notice that the closing position of the seed is currently unknown and will be determined only later, but the information in  $Z$  will suffice (together with some information available only when closing the interval) to reconstruct the basic encoding relative to the seed  $S$ , *i.e.*, the  $(\epsilon, S)$ -encoding.

In the second case,  $p$  is a closing position (lines 7–13) and the procedure removes from the stack  $Z$  all the records  $([b, e_1), l_S, b_\$)$  corresponding to an  $S$ -interval  $[b, p)$  whose

---

**Algorithm 1:** BuildBasicArcIntervals

---

**Input** : Three lists  $\mathcal{B}$ ,  $\mathcal{L}$ , and  $\mathcal{SA}$  containing the BWT, the LCP array and the GSA of the set  $R$ , respectively.

**Output**: A set of lists  $\mathcal{E}(\sigma, l_S, l_S)$  each containing the  $(\epsilon, S)$ -encoding for the seeds  $S$  whose first character is  $\sigma$  and whose length is  $l_S$ . The encodings  $([b, e), \cdot, \cdot, \cdot, \cdot, \cdot)$  in each list are sorted by increasing values of  $b$ .

```
1   $\#_{\$} \leftarrow 0$ ;
2  if  $\mathcal{B}[1] = \$$  then
3     $\#_{\$} \leftarrow 1$ ;
4   $p \leftarrow 2$ ;
5   $Z \leftarrow$  empty stack;
6  while  $p \leq |\mathcal{L}|$  do
7    if  $\mathcal{L}[p] < \mathcal{L}[p-1]$  then
8       $([b, e_1), l_S, b_{\$}) \leftarrow \text{top}(Z)$ ;
9      while  $Z$  is not empty and  $l_S > \mathcal{L}[p]$  do
10       if  $\#_{\$} > b_{\$}$  then
11         append  $([b, e_1), [b_{\$}, \#_{\$}), [1, |\mathcal{B}| + 1), [1, |\mathcal{B}| + 1), 0, l_S)$  to the list
12          $\mathcal{E}(C^{-1}(p), l_S, l_S)$ ;
13       pop( $Z$ );
14        $([b, e_1), l_S, b_{\$}) \leftarrow \text{top}(Z)$ ;
15    if  $\mathcal{L}[p] > \mathcal{L}[p-1]$  then
16       $(k, j) \leftarrow \mathcal{SA}[p-1]$ ;
17       $(k^*, j^*) \leftarrow \mathcal{SA}[p]$ ;
18       $q \leftarrow p$ ;
19      if  $\mathcal{L}[p] = k$  then
20        while  $\mathcal{L}[q+1] = \mathcal{L}[p] = k^* = k$  do
21           $q \leftarrow q + 1$ ;
22           $(k^*, j^*) \leftarrow \mathcal{SA}[q]$ ;
23        push  $([p-1, q+1), \mathcal{L}[p], \#_{\$})$  to  $Z$ ;
24       $p \leftarrow q$ ;
25    if  $\mathcal{B}[p] = \$$  then
26       $\#_{\$} \leftarrow \#_{\$} + 1$ ;
27       $p \leftarrow p + 1$ ;
28  while  $Z$  is not empty do
29     $([b, e_1), l_S, b_{\$}) \leftarrow \text{top}(Z)$ ;
30    if  $\#_{\$} > b_{\$}$  then
31      append  $([b, e_1), [b_{\$}, \#_{\$}), [1, |\mathcal{B}| + 1), [1, |\mathcal{B}| + 1), 0, l_S)$  to the list
32       $\mathcal{E}(C^{-1}(p), l_S, l_S)$ ;
33    pop( $Z$ );
```

---

opening position is  $b$  and whose forward  $\$$ -extension is  $[b, e_1)$ . The backward  $\$$ -extension is  $[b_{\$}, \#_{\$})$ , since  $b_{\$}$  is the number of sentinels in  $B[: b-1]$ , while  $\#_{\$}$  is the number of

sentinels in  $B[: p - 1]$ . Clearly  $[b_{\$}, \#_{\$}]$  is nonempty (and  $S$  is a seed) if and only if  $\#_{\$} > b_{\$}$ .

Notice that the stack  $Z$  always contains a nested hierarchy of distinct seeds (whose ending position might be currently unknown), that all  $k$ -intervals whose closing position is  $p$  are exactly the intervals with  $k > \mathcal{L}[p]$ , and they are found at the top of  $Z$ .

After all iterations, the stack  $Z$  contains the intervals whose closing position is  $p = |\mathcal{L}|$ . Those intervals are managed at lines 27–31.

There is a final technical detail: each output basic encoding associated to the overlap  $S$  is output to the list  $\mathcal{E}(S[1], |S|, |S|)$ . In fact we will use some different lists  $\mathcal{E}(\sigma, l_S, l_{PS})$ , each containing the encodings corresponding to seed  $S$  and left extension  $P$ , where  $\sigma$  is the first character of  $PS$ ,  $l_S = |S|$  and  $l_{PS} = |P| + |S|$ .

Moreover, a list  $\mathcal{E}(\sigma, l_S, l_{PS})$  is *correct* if it contains exactly the  $(P, S)$ -encodings such that  $|S| = l_S$  and  $|S| + |P| = l_{PS}$  and the encodings  $([b, e), \cdot, \cdot, \cdot, \cdot)$  are sorted by increasing values of  $b$ .

**Lemma 7.** *Let  $R$  be a set of reads, and let  $S$  be a seed of  $R$ . Then the  $(\epsilon, S)$ -encoding  $(q(S\$), q(\$S), [1, |\mathcal{B}| + 1), [1, |\mathcal{B}| + 1), 0, |S|)$  is output by Algorithm 1.*

*Proof.* Let  $[b_S, e_S)$  be the  $S$ -interval, let  $[b_{S\$}, e_{S\$})$  be the  $S\$$ -interval, and let  $[b_{\$S}, e_{\$S})$  be the  $\$S$ -interval. Since  $S$  is a seed, all those intervals are nonempty. Moreover, the sentinel is the smallest character, hence  $b_{S\$} = b_S$ .

When  $p = b_S + 1$ , since  $S$  is not an entire read (by definition of seed),  $S$  is a prefix of both  $LS[b_S + 1]$  and  $LS[b_S]$ , hence  $\mathcal{L}[p] \geq |S|$ . Moreover, since the  $S\$$ -interval is not empty,  $S\$$  is a prefix of  $LS[b_S]$  hence  $\mathcal{L}[p] = |S|$ , as the sentinel is not part of a common prefix.

By definition of  $S$ -interval,  $\mathcal{L}[p - 1] < |S|$ , hence the condition at line 14 is satisfied and at line 15  $k = \mathcal{L}[p] = |S|$ . When reaching line 22,  $\mathcal{L}[x] = k = |S|$  for all  $x$  with  $p \leq x \leq q$ . All those facts and the observation that  $S\$$  is a prefix of  $LS[p - 1]$ , imply that  $S\$$  is a prefix of all suffixes  $LS[i]$  with  $p - 1 \leq i \leq q$ . Since the while condition does not currently hold (as we have exited from the while loop),  $S\$$  is not a prefix of  $LS[q + 1]$ , hence  $[p - 1, q + 1)$  is the  $S\$$ -interval. Consequently at line 14 we push the triple  $(q(S\$), |S|, \#_{\$})$  on  $Z$ , where  $\#_{\$}$  is the number of sentinels in  $\mathcal{B}[: p]$ .

We distinguish two cases: either  $e_S \leq n$  or  $e_S > n$ . If  $e_S \leq n$  then there is an iteration where  $p = e_S$ . During such iteration the condition at line 7 holds, hence all entries  $(q(T\$), |T|, T_{\$})$  at the top of  $Z$  such that  $|T| > \mathcal{L}[p]$  are popped and the corresponding encoding is output at line 11. Since  $[b_S, e_S)$  is an  $S$ -interval,  $|S| > \mathcal{L}[p]$  hence the interval  $(q(S\$), |S|, \#_{\$})$  is popped. Since  $\$$  is the first symbol of the alphabet,  $q(\$S) = [b_{S\$}, e_{S\$}) = [b_{\$}, \#_{\$})$ .

If  $e_{S\$} > |\mathcal{L}|$ , then the condition of line 7 is never satisfied. Anyway, the stack  $Z$  is completely emptied at lines 27–31 and the same reasoning applies to show that the  $(\epsilon, S)$ -encoding is output.  $\square$

**Lemma 8.** *Let  $R$  be a set of reads, and let  $([b, e_1), [b_{\$}, \#_{\$}), [1, |\mathcal{B}| + 1), [1, |\mathcal{B}| + 1), 0, l_S)$  be an encoding output by Algorithm 1. Then  $q(S\$) = [b, e_1)$ ,  $q(\$S) = [b_{\$}, \#_{\$})$ ,  $l_S = |S|$  for some seed  $S$  of  $R$ .*

*Proof.* Notice that encodings are output only if a triple  $([p-1, q+1], \mathcal{L}[p], \#_{\$})$  is pushed on  $Z$ , which can happen only if  $\mathcal{L}[p] > \mathcal{L}[p-1]$ . By Lemma 5,  $p-1$  can be the opening position only of the seed  $S$  obtained by taking the  $\mathcal{L}[p]$ -long prefix of  $LS[p]$ . By the condition at line 18, the triple  $([p-1, q+1], \mathcal{L}[p], \#_{\$})$  is pushed on  $Z$  only if the  $S\$$ -interval is not empty.

Since  $\#_{\$} > b_{\$}$  and by the value of  $\#_{\$}$ , also the  $\$S$ -interval is nonempty, hence  $S$  is a seed.  $\square$

**Lemma 9.** *Let  $\mathcal{E}(\sigma, l_{PS}, l_{PS})$  be a list output by Algorithm 1. Let  $f_1 = (q(S_1\$), q(\$S_1), \cdot, \cdot, \cdot, |S_1|)$  and  $f_2 = (q(S_2\$), q(\$S_2), \cdot, \cdot, \cdot, |S_2|)$  be two encodings in  $\mathcal{E}(\sigma, 0, l_{PS})$ ,  $q(S_1\$) = [b_1, e_1)$  and  $q(S_2\$) = [b_2, e_2)$ . Then the intervals  $q(S_1\$)$  and  $q(S_2\$)$  are disjoint. Moreover,  $f_1$  precedes  $f_2$  in  $\mathcal{E}(\sigma, l_S, l_{PS})$  iff  $b_1 < b_2$ .*

*Proof.* By construction,  $l_S = |S_1| = |S_2| = l_{PS}$  and  $\sigma = S_1[1] = S_2[1]$ . Since  $|S_1| = |S_2|$ , the two string-intervals  $q(S_1\$)$  and  $q(S_2\$)$  cannot be nested, hence they are disjoint.

Notice that, since  $[b_1, e_1)$  and  $[b_2, e_2)$  are disjoint, then  $b_1 < e_1 \leq b_2$  or  $b_1 \geq e_2 > b_2$ . Assume that  $b_1 < e_1 \leq b_2$  and let us consider the iteration when  $p = e_1$ , *i.e.*, when  $f_1$  is output. Since  $e_1 \geq b_2$ , the entry  $([b_2, e_2 - 1], \cdot, \cdot, \cdot)$  has not been pushed to  $Z$  yet, hence  $f_1$  precedes  $f_2$  in  $\mathcal{E}(\sigma, 0, l_{PS})$ .

If  $b_1 \geq e_2 > b_2$  the same argument shows that  $f_2$  precedes  $f_1$  in  $\mathcal{E}(\sigma, 0, l_{PS})$ , completing the proof.  $\square$

**Corollary 10.** *Let  $\mathcal{E}(\sigma, 0, l_{PS})$  be a list computed by Algorithm 1. Then  $\mathcal{E}(\sigma, 0, l_{PS})$  contains exactly the  $(\epsilon, S)$ -encodings of all seed  $S$  such that  $\sigma = S[1]$ .*

There is an important observation on the sorted lists of encoding that we will manage. Let  $f_1 = (q(P_1S_1\$), q(\$S_1), q(P_1), q^r(P_1^r), |P_1|, |S_1|)$  and  $f_2 = (q(P_2S_2\$), q(\$S_2), q(P_2), q^r(P_2^r), |P_2|, |S_2|)$  be two encodings that are stored in the same list  $\mathcal{E}(\sigma, l_S, l_{PS})$ , hence  $|S_1| = |S_2|$  and  $|P_1S_1| = |P_2S_2|$ . Since  $|P_1S_1| = |P_2S_2|$ , the two string-intervals  $q(P_1S_1\$)$  and  $q(P_2S_2\$)$  are disjoint (as long as we can guarantee that  $P_1S_1 \neq P_2S_2$ ), hence sorting them by their opening position of the interval implies sorting also by closing position.

### 3.1 Labeling the overlap graph.

To complete the encoding of each arc, we need to compute the left extension. Such step will be achieved with the **ExtendEncodings** procedure (Algorithm 2), where the  $(P, S)$ -encodings are elaborated, mainly via backward  $\sigma$ -extensions, to obtain the  $(\sigma P, S)$ -encodings. Moreover, when  $PS$  has a nonempty backward  $\$$ -extension, we have determined that the encoding is terminal, hence we output the arc encoding to the lists  $\mathcal{A}(|P|, z)$  which will contain the arc encodings of the arcs incoming in the read  $r_z$  and whose left extension has length  $|P|$ .

The first fundamental observation is that a  $(P, S)$ -encoding can be obtained by extending a  $(\epsilon, S)$ -encoding with (if  $|P| = 1$ ), or by extending a  $(P_1, S)$ -encoding (if  $|P| > 1$  and  $P_1 = P[2:]$ ). Those extensions are computed in two phases: the first phase computes

---

**Algorithm 2:** ExtendEncodings( $l_P$ )

---

**Input** : Two lists  $\mathcal{B}$  and  $\mathcal{SA}$  containing the BWT and the GSA of the set  $R$ , respectively. The correct lists  $\mathcal{E}(\cdot, \cdot, \cdot)$  containing the  $(P, S)$ -encodings such that  $|P| = l_P$ .

**Output**: The correct lists  $\mathcal{P}(\cdot, \cdot, \cdot)$  containing the partially extended  $(\sigma P, S)$ -encoding. The arcs of the overlap graph outgoing from reads of length  $l_{PS} - 1$ , incoming in a read  $r_z$ , and with left extension long  $l_P$  are appended to the file  $\mathcal{A}(l_P, z)$ .

```
1  $\Pi, \pi \leftarrow |\Sigma|$ -long vectors  $\bar{0}$ ;
2  $p \leftarrow 0$ ;
3 foreach  $([b, e), q(\$S), q(P), q^r(P^r), l_P, l_{PS}) \in$   
    $\text{Merge}(\{\mathcal{E}(\sigma, l_{PS} - l_P, l_{PS}) : \sigma \in \Sigma, l_{PS} \geq l_P\})$  do
4    $\Pi(\sigma) \leftarrow \Pi[\sigma] + \pi[\sigma]$ , for each  $\sigma \in \Sigma$ ;
5   while  $p < b$  do
6      $\Pi[B[p]] \leftarrow \Pi[B[p]] + 1$ ;
7      $p \leftarrow p + 1$ 
8    $\pi \leftarrow \bar{0}$ ;
9   for  $p \leftarrow b$  to  $e - 1$  do
10    if  $\mathcal{B}[p] \neq \$$  then
11      Increment  $\pi[\mathcal{B}[p]]$  by 1;
12    if  $\mathcal{B}[p] = \$$  and  $l_P > 0$  then
13       $(k, j) \leftarrow \mathcal{SA}[p]$ ;
14      foreach read  $r_z \in q(\$S)$  do
15        Append the arc  $\langle j, i, q^r(P^r) \rangle$  to  $\mathcal{A}(l_P, z)$ ;
16    foreach  $\sigma \in \Sigma$  such that  $\pi[\sigma] > 0$  do
17      if  $l_P = l_{PS}$  then
18         $q(P) \leftarrow q(\sigma)$ ;  $q^r(P^r) \leftarrow q(\sigma)$ 
19       $b' \leftarrow C[\sigma] + \Pi[\sigma] + 1$ ;
20       $e' \leftarrow b' + \pi[\sigma]$ ;
21      Append  $([b', e'), q(S\$), q(P), q^r(P^r), l_P + 1, l_{PS} + 1)$  to  
       $\mathcal{P}(C^{-1}(b'), l_{PS} - l_P, l_{PS} + 1)$ ;
```

---

all partially extended encodings  $(q(\sigma PS\$), q(\$S), q(P), q^r(P^r), |P|, |S|)$  of the  $(P, S)$ -encodings. The second phase starts from the partially extended encodings and completes the extensions obtaining all  $(\sigma P, S)$ -encodings.

We iterate the procedure **ExtendEncodings** for increasing values of  $l_{PS}$ , where each step scans all lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$ , and writes the lists  $\mathcal{P}(\cdot, \cdot, l_{PS} + 1)$ , by computing all backward  $\sigma$ -extensions. The lists are called  $\mathcal{P}$  as a mnemonic for the fact that those encodings have been extended only partially. Those lists will be then fed to the **CompleteExtensions** procedure to complete the extensions, storing the result in the lists  $\mathcal{E}(\cdot, \cdot, l_{PS} + 1)$ .

The procedure **ExtendEncodings** basically extends a sequence of  $PS\$$ -intervals  $[b, e)$  that are sorted by increasing values of  $b$ . The procedure consists of a few parts; up

to line 11 the procedure maintains the arrays  $\Pi$  and  $\pi$  that are respectively equal to the number of occurrences of each symbol  $\sigma$  in  $\mathcal{B}[b : b - 1]$  (resp. in  $\mathcal{B}[b : p - 1]$ ). The correctness of this part is established by Lemmas 12, 11.

Lines 12–15 determine if the representative of  $[b, e)$  corresponds to an entire read, *i.e.*, if the current encoding is terminal; in that case some arcs of the overlap graph have been found and are output to the appropriate list.

The third part (lines 16–21) computes all backward  $\sigma$ -extensions of the current  $PS$ -interval  $[b, e)$ , obtaining the partially extended encodings. At line 21 we call the procedure that completes the extensions of the encodings.

In the following we will say that a list  $\mathcal{E}(\sigma, l_S, l_{PS})$  of encodings is correct if it contains exactly all  $(P, S)$ -encodings such that  $\sigma$  is the first character of  $PS$ ,  $l_P = |P|$ , and  $l_{PS} = |PS|$ . Moreover the encodings  $([b, e), \cdot, \cdot, \cdot, \cdot)$  are sorted by increasing values of  $b$ .

A list  $\mathcal{P}(\sigma, l_S, l_{PS})$  of partially extended encodings is correct if it contains exactly all partially extended  $(\sigma P, S)$ -encodings such that  $l_P = |P|$ , and  $l_{PS} = |PS|$ . Moreover the partially extended encodings  $([b, e), \cdot, \cdot, \cdot, \cdot)$  are sorted by increasing values of  $b$ .

Finally, we would like to point out that each list  $\mathcal{E}(\cdot, \cdot, \cdot)$  and  $\mathcal{P}(\cdot, \cdot, \cdot)$  contains disjoint intervals. If we can guarantee that the intervals in each list are sorted in non-decreasing order of the end boundary (we will prove this property of **CompleteExtensions**), then those intervals are also sorted in non-decreasing order of the start boundary (as required for the correctness of successive iterations of **ExtendEncodings**).

**Lemma 11.** *Let the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$  be the input of Algorithm 2 and assume that all those lists are correct. Let  $([b, e) = q(PS\$), q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  be the current encoding. Then at line 16 of Algorithm 2,  $\pi[\sigma]$  is equal to the number of occurrences of  $\sigma$  in  $\mathcal{B}[b : e - 1]$ .*

*Proof.* Notice that  $\pi$  is reset to zero at line 8 and is incremented only at line 11. The condition of the for loop (line 9) implies the lemma.  $\square$

**Lemma 12.** *Let the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$  be the input of Algorithm 2 and assume that all those lists are correct. Let  $([b, e) = q(PS\$), q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  be the current encoding. Then at line 8 of Algorithm 2,  $\Pi[\sigma]$  is equal to the number of occurrences of  $\sigma$  in  $\mathcal{B}[b : b - 1]$ .*

*Proof.* We will prove the lemma by induction on the number of the encodings that have been read. When extending the first input encoding,  $\pi$  consists of zeroes and the lemma holds, since the while loop at lines 5–7 increments  $\Pi$  by the number of occurrences of each symbol  $\sigma$  up to  $b - 1$ .

Let  $k$  be the number of encodings that have been read, with  $k \geq 2$ , and let  $([b_1, e_1), \cdot, \cdot, \cdot, \cdot)$  be the  $(k - 1)$ -th encoding read. By Lemma 11,  $\pi[\sigma]$  is equal to the number of occurrences of  $\sigma$  in  $\mathcal{B}[b_1, e_1)$ , hence after line 4  $\Pi$  contains the number of occurrences of each symbol in  $\mathcal{B}[b : p - 1]$ . The while loop at lines 5–7 increments  $\Pi$  by the number of occurrences of each symbol  $\sigma$  in the portion of  $\mathcal{B}$  between  $e_1$  and  $b - 1$ , completing the proof.  $\square$

**Lemma 13.** *Let the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$  be the input of Algorithm 2 and assume that all those lists are correct. Let  $([b, e) = q(PS\$), q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  a generic*

input encoding, and let  $\sigma$  be a character of  $\Sigma$ . If  $[b, e]$  has a nonempty backward  $\sigma$ -extension  $[b_1, e_1]$ , then **ExtendEncodings** outputs the partially extended encoding  $([b_1, e_1], q(\$S), q(P), q^r(P^r), l_P, l_{PS} + 1)$  to the list  $\mathcal{P}(\sigma, |S|, l_{PS+1})$ .

*Proof.* By Lemma 3.1 in [17], the backward  $\sigma$ -extension of  $[b, e]$  is equal to  $[C[\sigma] + \text{Occ}(\sigma, b) + 1, C[\sigma] + \text{Occ}(\sigma, e)]$ . By Lemmas 12, 11, the values of  $b'$  and  $e'$  computed at lines 19–20 is correct.

Notice that  $\pi[\sigma] > 0$  iff  $\text{Occ}(\sigma, e) > \text{Occ}(\sigma, b) + 1$ , hence the partially extended encoding  $([b_1, e_1], q(\$S), q(P), q^r(P^r), l_P, l_{PS} + 1)$  is output iff  $e_1 > b_1$ , that is the backward  $\sigma$ -extension is nonempty.  $\square$

**Lemma 14.** Let  $([b_1, e_1], q(\$S), q(P), q^r(P^r), l_P, l_{PS} + 1)$  be a partially extended encoding that is output by **ExtendEncodings**. Let  $\sigma = C^{-1}(b_1)$ . Then  $([b, e], q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  is an encoding in  $\mathcal{P}(\sigma, |S|, l_{PS})$  and  $[b_1, e_1]$  is the backward  $\sigma$ -extension of  $[b, e]$ .

*Proof.* Let  $([b, e], q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  be the current input encoding when  $([b_1, e_1], q(\$S), q(P), q^r(P^r), l_P, l_{PS} + 1)$  is output by **ExtendEncodings**. When computing the partially extended interval (line 21),  $b_1 = C[\sigma] + \Pi[\sigma] + 1$  and  $e_1 = b_1 + \pi[\sigma]$ . The lemma is a direct consequence of Proposition 3 and Lemmas 11, 12.  $\square$

**Lemma 15.** Let  $\mathcal{P}(\sigma, l_{PS} - l_P, l_{PS} + 1)$  be any list written by **ExtendEncodings**. Then the encodings  $([b', e'], q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  in  $\mathcal{E}(\sigma, l_{PS} - l_P, l_{PS} + 1)$  are sorted by increasing values of  $e'$ . Moreover the intervals  $[b', e']$  are disjoint.

*Proof.* By Lemmas 13 15, the list  $\mathcal{P}(\sigma, l_{PS} - l_P, l_{PS} + 1)$  contains only partially extended encodings relative to the pairs  $(P, S)$  where  $\sigma$  is the first symbol of  $PS$ .

Let us now consider a generic partially extended encoding  $([b', e'], q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  written to  $\mathcal{P}(\sigma, l_{PS} - l_P, l_{PS} + 1)$ . Notice that  $e' - b' = \pi[\sigma]$  and that, while managing the next partially extended encoding,  $\Pi[\sigma]$  will be incremented by  $e' - b'$ , hence during the next iterations of the foreach loop  $C[\sigma] + \Pi[\sigma] + 1$  will be at least as large as the value of  $e'$  at the current iteration. That completes the proof, since the start boundary is always equal to  $C[\sigma] + \Pi[\sigma] + 1$  (see line 19).  $\square$

The following corollary summarizes this subsection.

**Corollary 16.** Let the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$  be the input of Algorithm 2 and assume that all those lists are correct. Let  $([b, e], q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  be the current encoding.

Then **ExtendEncodings** produces the correct lists  $\mathcal{P}(\sigma, l_s, l_{PS} + 1)$ .

**Lemma 17.** If the input encodings are correct, then Algorithm 2 outputs the arc  $(r_j, r_i)$  to  $\mathcal{A}(l_P, i)$  iff there exists a read  $r_j = PS$  with  $P$  and  $S$  both nonempty and there exists a read  $r_i$  whose prefix is  $S$ .

*Proof.* Notice that the encoding of the arc  $(r_j, r_i)$  is output to  $\mathcal{A}(l_{PS} + 1, i)$  only if we are currently extending the  $(P, S)$ -encoding (hence  $S$  is a seed) and we have found that the  $PS$ -interval has a nonempty backward  $\sigma$ -extension, since  $B[p]$  is the symbol preceding

$PS$  in a suffix and  $B[p] = \$$ . By definition of seed,  $S$  is nonempty and there exists a read  $r_i$  with prefix  $S$ , while  $P$  is nonempty by the condition at line 12.

Assume now that  $r_j = PS$  is a read with  $P$  and  $r_i$  is a read with prefix  $S$ . The  $S\$$ -interval and the  $PS\$$ -interval are nonempty. Moreover  $R^p(S) \neq \emptyset$  and  $S$  a seed, hence there is an iteration of **ExtendEncodings** where we backward extend the  $PS\$$ -encoding. Since  $r_j = PS$ ,  $PS$  has a nonempty backward  $\$$ -extension, hence the condition at line 12 is satisfied.  $\square$

### 3.2 Extending arc labels

While the procedure **ExtendEncodings** backward extends the  $PS\$$ -intervals, the actual arc labels are the  $P$ -intervals, therefore we need a dedicated procedure, called **CompleteExtensions**, that scans the results of **ExtendEncodings**, *i.e.*, a list of partially extended encodings and updates them by extending the intervals  $q(P)$  on  $R$  and  $q^r(P^r)$  on  $R^r$ .

A procedure **ExtendIntervals** has been originally described [6] to compute all backward extensions of a set of *disjoint* string-intervals, with only a single pass over  $\mathcal{B}$ . In our case, the string-intervals are not necessarily disjoint, therefore that procedure is not directly applicable. We exploit the property that any two string-intervals are either nested or disjoint to design a new procedure that computes all backward extensions with a single scan of the list  $\mathcal{B}$ .

Our procedure **CompleteExtensions** takes in input a list  $I$  of partially extended encodings  $([b, e), q(S\$), q(P), q^r(P^r), l, l_{PS})$ , sorted by increasing values of  $b$  and (as a secondary criterion) by decreasing values of  $e$ . Moreover the list  $I$  is terminated with a sentinel partially extended encoding  $(\cdot, \cdot, [n+1, n+2), \cdot, \cdot, \cdot)$ —we recall that  $n$  is the total number of input characters. For each input encoding coming from the list  $\mathcal{P}(\sigma, \cdot, \cdot)$ , the procedure outputs the record  $(p(PS\$), q(S\$), q(\sigma P), q^r(P^r \sigma), l_P, l_{PS})$ .

Just as the procedure **ExtendEncodings**, we maintain an array  $\Pi$ , where  $\Pi[\sigma]$  is equal to the number of occurrences of the character  $\sigma$  in  $\mathcal{B}[p-1]$ , and  $p$  is the current position in  $\mathcal{B}$ . The only difference w.r.t. **ExtendEncodings** is that  $\sigma$  can be the sentinel character  $\$$ . We recall that  $Occ(\sigma, p)$  is the number of occurrences of  $\sigma$  in  $\mathcal{B}[p-1]$  [17], therefore  $\Pi[\sigma] = Occ(\sigma, p)$ , where  $p$  is the number of symbols of  $\mathcal{B}$  that have been read. The array  $\Pi$  is used to compute the backward extension at line 16 of Algorithm 3.

We also maintain a stack  $Z$  storing the partially extended encodings that have already been read, but have not been extended yet. A correct management of  $Z$  allows to have the encodings in the correct order, that is to read the encodings in increasing order of  $b$ , and to actually extend the encodings in increasing order of  $e$ . This ordering allows to scan sequentially  $\mathcal{B}$ .

We will start by showing that **CompleteExtensions** correctly manages the array  $\Pi$ .

**Lemma 18.** *Assume that the input partial encodings are  $([b, e) = q(\sigma PS\$), q(S\$), q(P), q^r(P^r), l_P, l_{PS})$  and that they are sorted by increasing value of  $e$ . Then after lines 11 and 19, and before lines 9 and 17 of Algorithm 3,  $\Pi[\sigma] = Occ(\sigma, p)$ .*



---

**Algorithm 3:** CompleteExtensions( $l_P$ )

---

**Input** : The BWT  $B$  of a set  $R$  of strings. The correct lists  $\mathcal{P}(\cdot, \cdot, \cdot)$  containing all partially extended  $(\sigma P, S)$ -encoding such that  $|P| = l_P$ .

**Output** : The correct lists  $\mathcal{E}(\cdot, \cdot, \cdot)$  containing all  $(\sigma P, S)$ -encodings.

```
1  $I \leftarrow \text{Merge}(\{\mathcal{P}(\sigma, l_{PS} - l_P, l_{PS} + 1) : \sigma \in \Sigma, l_{PS} \geq l_P\})$ ;  
2 Append the sentinel interval  $(\cdot, \cdot, [n + 1, n + 2), \cdot, \cdot, \cdot)$  to  $I$ ;  
3  $\Pi \leftarrow |\Sigma|$ -long vector  $\bar{0}$ ;  
4  $Z \leftarrow$  stack with the record  $\langle(\cdot, \cdot, [1, \infty), \cdot, \cdot, \cdot), \Pi\rangle$ ;  
5  $p \leftarrow 1$ ;  $e_z \leftarrow +\infty$ ;  
6 foreach  $(q_1, q_2, [b, e), [b', e^r), l_P, l_{PS}) \in I$  do  
7    $\langle(\cdot, \cdot, [e_z), \cdot, \cdot, \cdot), \cdot\rangle \leftarrow \text{top}(Z)$ ;  
8   while  $e > e_z$  do  
9     while  $p < e_z - 1$  do  
10       $\Pi[B[p]] \leftarrow \Pi[B[p]] + 1$   $p \leftarrow p + 1$ ;  
11       $p \leftarrow e_z$ ;  
12       $\langle(q_{1z} = [b_{ps}, e_{ps}), q_{2z}, [b_z, e_z), [b'_z, e_z^r), l_{pz}, l_{psz}), \Pi_z\rangle \leftarrow \text{pop}(Z)$ ;  
13       $\sigma \leftarrow C^{-1}(b_{ps})$ ;  
14       $prev \leftarrow \sum_{c < \sigma} (\Pi(c) - \Pi_z(c))$ ;  
15       $w \leftarrow \Pi(\sigma) - \Pi_z(\sigma)$ ;  
16      Append  $(q_{1z}, q_{2z}, [C[\sigma] + \Pi_z[\sigma] + 1, C[\sigma] + \Pi_z[\sigma] + 1 + w), [b'_z + prev, b'_z + prev + w), l_{pz}, l_{ez})$  to the list  $\mathcal{E}(\sigma, l_{psz} - l_{pz}, l_{psz})$ ;  
17     while  $p < b$  do  
18        $\Pi[B[p]] \leftarrow \Pi[B[p]] + 1$ ;  
19        $p \leftarrow p + 1$ ;  
20   push( $Z, \langle(q_1, q_2, [b, e), [b', e^r), l_P, l_{PS}), \Pi\rangle$ );
```

---

*Proof.* We can prove the lemma by induction on the number  $i$  of input encodings that have been read so far. Notice that the lemma holds at line 5, since  $\Pi = \bar{0}$  and  $p = 1$ . When  $i = 1$ , the condition at line 8 does not hold, hence we only need to consider the value of  $\Pi$  at lines 17 and 19. A direct inspection of lines 17–19 shows that the lemma holds in this case.

Assume now that  $i > 1$ . Since the procedure modifies  $\Pi$  or  $p$  only at lines 9–11 and lines 17–19, by inductive hypothesis the lemma holds before line 9. Again, a direct inspection of lines 9–11 proves that the lemma holds at line 11, which in turn implies that it holds at line 17. The same direct inspection of lines 17–19 as before is sufficient to complete the proof.  $\square$

Since elements are pushed on the stack  $Z$  only at lines 4 and 20, and partially extended encodings are pushed on  $Z$  without any change, a direct consequence of Lemma 18 is the following corollary.

**Corollary 19.** *Let  $\langle(\cdot, \cdot, [b, e), \cdot, \cdot, \cdot), \Pi\rangle$  be an element of  $Z$ , and let  $\sigma$  be any character in  $\Sigma^\$$ . Then  $\Pi[\sigma] = \text{Occ}(\sigma, b)$ .*

---

**Algorithm 4:** OverlapGraph( $R$ )

---

**Input** : A set  $R$  of reads.  
**Output**: The overlap graph  $G_O$  of  $R$ .  
1  $l_{max} \leftarrow$  the maximum length of a read in  $R$ ;  
2 Construct  $\mathcal{B}, \mathcal{L}, \mathcal{SA}$ ;  
3 BuildBasicArcIntervals( $R$ ); /\* computes  $\mathcal{E}(\cdot, l_{PS}, l_{PS})$  \*/  
4 **for**  $l_P \leftarrow 0$  **to**  $l_{max} - 2$  **do**  
5     ExtendEncodings( $l_P$ ); /\* computes  $\mathcal{P}(\cdot, l_{PS} - l_P, l_{PS})$  \*/  
6     CompleteExtensions( $l_P$ ); /\* computes  $\mathcal{E}(\cdot, l_{PS} - l_P, l_{PS})$  \*/

---

**Lemma 20.** *Let  $\mathcal{P}(\cdot, \cdot, l_{PS})$  be the correct lists that are the input of Algorithm 3, and let  $(q(\sigma PS\$), q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  be a generic partially extended encoding in one of such lists. Then Algorithm 3 outputs the encoding  $(q(\sigma PS\$), q(\$S), q(\sigma P), q^r(P^r\sigma), l_P, l_{PS})$  if and only if  $(q(\sigma PS\$), q(\$S), q(P), q^r(P^r), l_P, l_{PS})$  is an input partial encoding.*

*Proof.* In Algorithm 3 each input partially extended encoding is pushed on the stack  $Z$  exactly once and extended exactly once. Hence we only need to prove that for each partially extended input encoding  $(q(\sigma PS\$), q(\$S), q(\sigma P), q^r(P^r\sigma), l_P, l_{PS})$ , the  $(\sigma P, S)$ -encoding is output.

Let  $[b, e]$  be equal to  $q(P)$ , let  $[b', e']$  be equal to  $q^r(P^r)$  and let  $\sigma$  be the symbol  $C^{-1}(b_{ps})$ . Notice that the output encoding is obtained by computing  $q(\sigma P)$  and  $q^r(P^r\sigma)$ . By Proposition 3,  $q(\sigma P)$  is equal to  $[C[\sigma] + Occ(\sigma, b) + 1, C[\sigma] + Occ(\sigma, e)]$ , which is correctly computed at line 16 (by Lemma 18 and Corollary 19).

Moreover,  $q^r(P^r\sigma)$  is equal to  $[b'_1, e'_1]$ , by Proposition 3, Lemma 18 and Corollary 19. Moreover, the encoding  $(q(PS\$), q(\$S), [b_1, e_1], [b'_1, e'_1], l_P, l_{PS})$  is output at line 16.  $\square$

Notice that the value of  $p$  never decreases, since its value is modified only by increments. This fact implies that  $\mathcal{B}$  is scanned sequentially. To complete the correctness of our algorithm, we need to show that the output follows the desired ordering. We will start with some lemmas showing the structure of the encodings stored in the stack  $Z$ .

**Lemma 21.** *The stack  $Z$  of Algorithm 3 contains a hierarchy of encodings  $(\cdot, \cdot, [b, e], \cdot, \cdot, \cdot)$  where all intervals  $[b, e]$  are nested, with the smallest at the top.*

*Proof.* We only have to prove that the lemma holds at line 20, since it is the only line where an encoding is pushed on a nonempty stack. Let  $Z$  be the stack just before the push, and let  $(\cdot, \cdot, [b_z, e_z], \cdot, \cdot, \cdot)$  be the encoding at the top of  $Z$ .

Clearly the lemma holds when  $Z$  contains only the sentinel encoding pushed at line 4, therefore assume that the top encoding of  $Z$  is an input encoding.

Notice that an encoding is pushed on  $Z$  without modification, before reading the next input encoding. Since the input encodings are sorted by increasing values of  $b$ , then  $b \geq b_z$ . To reach line 20, the condition at line 8 must be false, hence  $e \leq e_z$ . Consequently  $[b, e]$  is included in  $[b_z, e_z]$ .  $\square$

**Lemma 22.** *Algorithm 3 pops all input encodings  $(\cdot, \cdot, [b, e], \cdot, \cdot, \cdot)$  in nondecreasing order of  $e$ , but it does not pop the sentinel encodings.*

*Proof.* First, we will consider a single generic iteration of the while loop at lines 8–16. By Lemma 21 the intervals in  $Z$  are nested, therefore the intervals popped in a single iteration satisfy the lemma.

We can consider the intervals popped in different iterations. Let  $f_1 = (\cdot, \cdot, [b_1, e_1], \cdot, \cdot, \cdot)$  be the most recently popped encoding, and let  $f_z = (\cdot, \cdot, [b_z, e_z], \cdot, \cdot, \cdot)$  be a generic interval that has been popped from  $Z$  in a previous iteration; we will show that  $e_z \leq e_1$ . Moreover let  $f = (\cdot, \cdot, [b, e], \cdot, \cdot, \cdot)$  be the encoding read from  $I$  at the iteration when  $f_1$  has been popped from  $Z$ . By construction,  $f_z$  precedes  $f$  (which precedes  $f_1$ ) in  $I$ . Since the intervals in  $I$  are in non-decreasing order of the start boundary,  $b_z \leq b \leq b_1$ . Moreover, the condition at line 8 that determines when to pop an encoding, implies that  $e > e_z$ . All string-intervals in  $I$  are disjoint or nested, therefore  $e \leq b_1$  or  $e \geq e_1$ . If  $e \leq b_1$ , then  $e \leq e_1$  and, a fortiori since  $e_z < e$ , then  $e_z \leq e_1$ . Hence we only need to consider the case when  $[b, e]$  includes  $[b_1, e_1]$ , that is  $e \geq e_1$ . Now, let us consider the intervals  $[b, e]$  and  $[b_z, e_z]$ . Since  $e_z < e$  and  $b_z \leq b$ , those intervals cannot be nested, hence  $e_z \leq b$ . Since  $b \leq b_1$  and  $b_1 < e_1$ , then  $e_z \leq e_1$ .

Finally, we want to prove that all intervals in  $I$ , except for the sentinel intervals, are popped from  $Z$  (and backward extended). Just after reading from  $I$  the sentinel  $(\cdot, \cdot, [n+1, n+2], \cdot, \cdot, \cdot)$ , all intervals in  $Z$ , but not the starting sentinel, satisfy the condition at line 8, completing the proof.  $\square$

**Corollary 23.** *The lists  $\mathcal{E}(\sigma, l_P, l_{PS})$  are correct.*

*Proof.* It is a direct consequence of Corollary 16 and Lemma 22.  $\square$

**Corollary 24.** *Algorithm 4 correctly computes the arcs of the overlap graph  $G_O$ .*

*Proof.* It is a direct consequence of Corollary 10 (which shows the correctness of Algorithm 1 to compute the basic encodings), Corollary 16 (which shows the correctness of Algorithm 2 to compute the partial extensions of a set of encodings), Lemma 17 (which shows the correctness of Algorithm 3 to complete the extension of a set of partially extended encodings), and finally Corollary 25 which shows that the arcs of the overlap graph  $G_O$  are correctly output.  $\square$

Algorithm 1 scans only once the lists  $\mathcal{B}, \mathcal{SA}, \mathcal{L}$  (hence reading  $3n$  records, with  $n = |\mathcal{B}|$ ) and outputs at most  $n$  records, by Lemma 5.

Each execution of Algorithm 2 scans only once the lists  $\mathcal{B}, \mathcal{SA}, \mathcal{L}$  (hence reading  $3n$  records) as well as the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$  containing the  $(P, S)$ -encodings with  $|PS| = l_{PS}$ . Let us now consider the  $(P, S)$ -encodings that are read during a single execution of Algorithm 2, and notice that the corresponding  $(P, S)$ -intervals are disjoint, since the length of  $|PS|$  is always equal to  $l_{PS}$ . This fact implies that at most  $n$   $(P, S)$ -encodings are read and at most  $n$   $(P, S)$ -encodings are output.

The analysis of Algorithm 3 is similar to that of Algorithm 2. The only difference is that Algorithm 3 scans once the list  $\mathcal{B}$ , as well as the lists  $\mathcal{E}(\cdot, \cdot, l_{PS})$ . The consequence is that Algorithm 3 reads at most  $2n$  records and writes at most  $n$  records.

Since Algorithms 2 and Algorithm 3 are called at most  $l$  times, the overall number of records that are read is at most  $3n + 6ln$ . Notice that this I/O complexity matches the one of BCReXt [6], which is the most-efficient known external-memory algorithm to compute the data structures (GSA, BWT, LCP) we use to index the input reads.

**Corollary 25.** *Given the lists  $\mathcal{B}$ ,  $\mathcal{SA}$ ,  $\mathcal{L}$ , it is possible to compute the overlap graph of a set of reads  $R$  with total length  $n$  and where no read is longer than  $l$  characters, reading sequentially  $(3 + 6l)n$  records.*

## 4 Reducing the overlap graph to a string graph

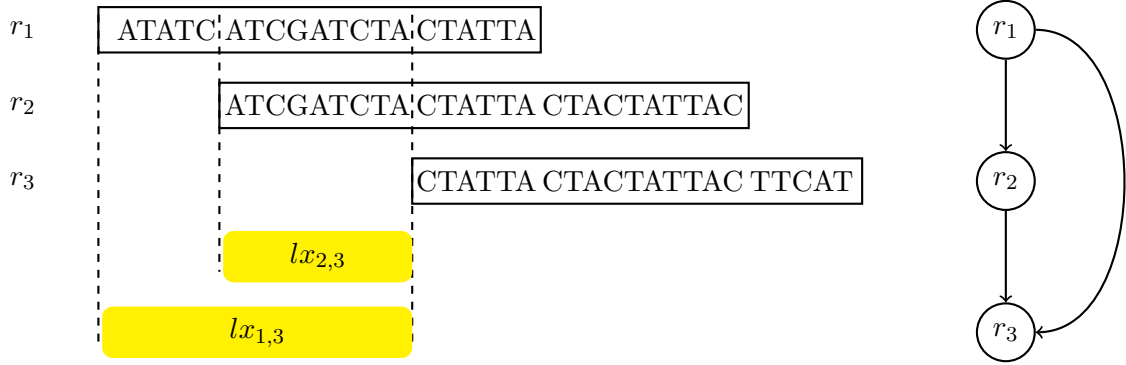


Figure 2: Example of reducible arc of the overlap graph. The read  $r_1$  is equal to ATATCATCGATCTACTATTA, while the read  $r_2$  is equal to ATCGATCTACTATTACTACTATTAC and the read  $r_3$  is CTATTACTACTATTACTTTCAT. The associated overlap graph is on the right. The arc  $(r_1, r_3)$  is reducible.

In this section we state a characterization of string graphs based on the notion of string-interval, then we will exploit such characterization to reduce the overlap graph.

**Lemma 26.** *Let  $G_O$  be the labeled overlap graph for a substring-free set  $R$  of reads and let  $(r_i, r_j)$  be an arc of  $G_O$ . Then,  $(r_i, r_j)$  is reducible iff there exists another arc  $(r_h, r_j)$  of  $G_O$  incoming in  $r_j$  and such that  $lx_{h,j}^r$  is a proper prefix of  $lx_{i,j}^r$ .*

*Proof.* First notice that  $lx_{h,j}^r$  is a proper prefix of  $lx_{i,j}^r$  iff and only if  $lx_{h,j}$  is a proper suffix of  $lx_{i,j}$ .

By definition,  $(r_i, r_j)$  is reducible if and only if there exists a second path  $r_i, r_{h_1}, \dots, r_{h_k}, r_j$  representing the string  $XYZ$ , where  $X$ ,  $Y$  and  $Z$  are respectively the left extension of  $r_j$  with  $r_i$ , the overlap of  $r_i$  and  $r_j$ , and the right extension of  $r_i$  with  $r_j$ .

Assume that such a path  $(r_i, r_{h_1}, \dots, r_{h_k}, r_j)$  exists. Since  $r_i, r_{h_1}, \dots, r_{h_k}, r_j$  represents  $XYZ$  and  $Z = rx_{i,j}$ ,  $r_{h_k} = X_1YZ_1$  where  $X_1$  is a suffix of  $X$  and  $Z_1$  is a proper prefix

of  $Z$ . Notice that  $X_1 = lx_{h_k,j}$  and  $R$  is substring free, hence  $X_1$  is a proper suffix of  $X$ , otherwise  $r_i$  would be a substring of  $r_{h_k}$ , completing this direction of the proof.

Assume now that there exists an arc  $(r_h, r_j)$  such that  $lx_{h,j}$  is a proper suffix of  $lx_{i,j}$ . Again,  $r_h = X_1 Y_1 Z_1$  where  $X_1$ ,  $Y_1$  and  $Z_1$  are respectively the left extension of  $r_j$  with  $r_h$ , the overlap of  $r_h$  and  $r_j$ , and the right extension of  $r_h$  with  $r_j$ . By hypothesis,  $X_1$  is a proper suffix of  $X$ . Since  $r_h$  is not a substring of  $r_i$ , the fact that  $X_1$  is a suffix of  $X$  implies that  $Y$  is a substring of  $Y_1$ , therefore  $r_i$  and  $r_h$  overlap and  $|ov_{i,h}| \geq |Y|$ , hence  $(r_i, r_h)$  is an arc of  $G_O$ .

The string associated to the path  $r_i, r_h, r_j$  is  $r_i r x_{i,h} r x_{h,j}$ . By Lemma 4,  $r_i r x_{i,h} r x_{h,j} = lx_{i,h} lx_{h,j} r_j$ . At the same time the string associated to the path  $r_i, r_j$  is  $r_i r x_{i,j} = lx_{i,j} r_j$  by Lemma 4, hence it suffices to prove that  $lx_{i,h} lx_{h,j} = lx_{i,j}$ . Since  $lx_{h,j}$  is a proper suffix of  $lx_{i,j}$ , by definition of left extension,  $lx_{i,h} lx_{h,j} = lx_{i,j}$ , completing the proof.  $\square$

Since the encoding of an arc  $(r_h, r_j)$  contains both  $q(lx_{h,j}^r)$  and  $|lx_{h,j}^r|$ , we can transform Lemma 26 into an easily testable property, by way of Proposition 2. The following Lemma 27 shows that, if  $(r_x, r_z)$  can be reduced, then it can be reduced by an arc of the string graph  $G_S$ , hence avoiding a comparison between all pairs of arcs of  $G_O$ .

**Lemma 27.** *Let  $G_O$  be the overlap graph of a string  $R$  of reads, let  $G_S$  be the corresponding string graph, and let  $(r_x, r_z)$  be an arc of  $G_O$  that is not an arc of  $G_S$ . Then there exists an arc  $(r_s, r_z)$  of  $G_S$  such that  $q^r(lx_{s,z}^r)$  includes  $q^r(lx_{x,z}^r)$  and  $|lx_{x,z}| > |lx_{s,z}|$ .*

*Proof.* Let  $(r_s, r_z)$  be the arc of  $G_O$  whose left extension is the shortest among all arcs of  $G_O$  such that  $q^r(lx_{s,z}^r)$  includes  $q^r(lx_{x,z}^r)$  and  $|lx_{x,z}| > |lx_{s,z}|$ . By Lemma 26, since  $(r_x, r_z)$  is not an arc of  $G_S$  such an arc must exist. We want to prove that  $(r_s, r_z)$  is an arc of  $G_S$ .

Assume to the contrary that  $(r_s, r_z)$  is not an arc of  $G_S$ , that is there exists an arc  $e_1 = (r_h, r_z)$  of  $G_O$  such that  $q^r(lx_{s,z})$  includes  $q^r(lx_{h,z})$  and  $|lx_{s,z}| > |lx_{h,z}|$ . Then  $q^r(lx_{x,z})$  includes  $q^r(lx_{h,z})$  and  $|lx_{x,z}| > |lx_{h,z}|$ , contradicting the assumption that  $(r_s, r_z)$  is the arc of  $G_O$  whose left extension is the shortest among all arcs of  $G_O$  such that  $q^r(lx_{s,z}^r)$  includes  $q^r(lx_{x,z}^r)$  and  $|lx_{x,z}| > |lx_{s,z}|$ .  $\square$

Lemma 27 suggests that each arc  $e$  of  $G_O$  should be tested only against arcs in  $G_S$  whose left extension is strictly shorter than that of  $e$  to determine whether  $e$  is also an arc of  $G_S$ . A simple comparison between each arc of the overlap graph and each arc of the string graph would determine which arcs are irreducible, but this approach would require to store in main memory all arcs of  $G_S$  incident on a vertex. To reduce the main memory usage, we partition the arcs of  $G_S$  incoming in a vertex  $r_z$  into chunks, where each chunk can contain at most  $M$  arcs (for any given constant  $M$ ) [32]. Let  $D_z$  be the set of arcs of  $G_S$  incoming in  $r_z$ , and let  $d_z$  be its cardinality. Since there are at most  $M$  arcs in each chunk, we need  $\lceil d_z/M \rceil$  passes over  $D_z$  to perform all comparisons. There are some technical details that are due to the fact that the set  $D_z$  is not known before examining the arcs of  $G_O$  incoming in  $r_z$  (see Algorithm 5). Mainly, we need an auxiliary file to store whether each arc  $e$  of  $G_O$  has already been processed, that is if we have already decided whether  $e$  is an arc of  $G_S$ .

Now we can start proving the the correctness of Algorithm 5.

---

**Algorithm 5:** ReduceOverlapGraph( $M$ )

---

**Input** : The number  $M$  of 5-integer records that can be stored in memory. The lists  $\mathcal{A}(p, \cdot)$  of arc encodings of  $G_O$ .

**Output** : The set  $E$  of arc encodings of the irreducible arcs of  $G_O$ .

```
1  $l_{\max}$  the maximum length of a read in  $R$ ;  
2  $|R|$  the number of reads in  $R$ ;  
3 for  $z \leftarrow 1$  to  $|R|$  do  
4    $D_z \leftarrow \emptyset$ ;  
5   while  $\exists p : \mathcal{A}(p, z)$  contains at least an encoding not marked processed do  
6      $C \leftarrow \emptyset$ ; /*  $C$  contains a chunk of edges of  $G_S$  */  
7     for  $p \leftarrow 1$  to  $l_{\max}$  do  
8       foreach unprocessed arc encoding  $e = \langle i, z, q^r(Q^r), |Q| \rangle \in \mathcal{A}(p, z)$  do  
9         foreach  $\langle h, z, q^r(P^r), |P| \rangle \in C$  do  
10          if  $|P| < |Q|$  and  $q^r(P^r)$  contains  $q^r(Q^r)$  then  
11            Mark  $e$  as transitive and processed;  
12          if  $|C| < M$  and  $e$  is not marked transitive then  
13            Add  $e$  to  $C$ ;  
14            Mark  $e$  as processed;  
15    $D_z \leftarrow D_z \cup C$ ;  
16 return  $\cup_z D_z$ ;
```

---

**Lemma 28.** Let  $G_O$  be the overlap graph of a string  $R$  of reads, and let  $G_S$  be the corresponding string graph. Then the execution of Algorithm 5 on  $G_O$  terminates with all arcs of  $G_O$  marked processed.

*Proof.* To prove that the algorithm terminates, we only have to prove that all arcs in the generic set  $\mathcal{A}(p, z)$  are marked processed, as in that case the condition of the while at line 5 becomes false. As long as there is an unprocessed arc, the condition at line 5 is satisfied, hence the corresponding while loop is executed. At each execution of such loop, the first unprocessed arc is added to  $C$  (since  $C$  is emptied at the beginning of the iteration) and marked processed. Hence, eventually all arcs must be marked processed.  $\square$

**Lemma 29.** Let  $G_O$  be the overlap graph of a string  $R$  of reads, let  $G_S$  be the corresponding string graph, and let  $e = \langle i, z, q^r(Q^r), |Q| \rangle$  be an arc encoding that is marked transitive. Then  $e$  is not an arcs of  $G_S$ .

*Proof.* Since  $e$  is marked transitive, there exists an arc encoding  $\langle h, z, q^r(P^r), |P| \rangle \in C$  such that  $|P| < |Q|$  and  $q^r(P^r)$  contains  $q^r(Q^r)$ . By construction of arc encoding and by Lemma 26, the arc  $(r_i, r_z)$  cannot be an arc of  $G_S$ .  $\square$

**Lemma 30.** Let  $G_O$  be the overlap graph of a string  $R$  of reads, let  $G_S$  be the corresponding string graph, and let  $e = \langle i, z, q^r(Q^r), |Q| \rangle$  be an arc encoding inserted into  $D_z$  by Algorithm 5. Then  $(r_i, r_z)$  is an arc of  $G_S$ .

*Proof.* Since  $e$  is in  $D_z$ , previously  $e$  has been added to  $C$ . Let us consider the iteration when  $e$  is added to  $C$ : notice that  $|C| < M$  and  $e$  is not marked as processed at the beginning of the iteration. Consequently, no arc encoding that is currently in  $C$  or that has been in  $C$  in a previous iterations of the while loop at lines 5–15 satisfies the condition of Lemma 26, that is no arc in  $C$  or in a previous occurrence of  $C$  can reduce  $e$ .

Since the arcs incoming in  $r_z$  are examined in increasing order of their left extension, all arcs of  $G_S$  that are incoming in  $r_z$  and whose left extension is shorter than  $e$  have already been inserted in  $C$ , either in the current iteration or in one of the previous iterations. Consequently no arc of  $G_O$  can reduce  $e$ , hence  $e$  is an arc of  $G_S$ .  $\square$

**Theorem 31.** *Let  $G_O$  be the overlap graph of a string  $R$  of reads. Then Algorithm 5 outputs the set  $E$  of the arc encodings of the irreducible arcs of  $G_O$  reading or writing at most  $3|E(G_O)|\lceil d/M \rceil$  records, where  $E(G_O)$  is the arc set of  $G_O$  and  $d$  is the maximum indegree of  $G_S$ .*

*Proof.* By Lemmas 28, 29, and 30, Algorithm 5 outputs the set  $E$  of the arc encodings of the irreducible arcs of  $G_O$ .

To determine the total number of records that are read by Algorithm 5, we notice that each execution of the while loop at lines 5–15 read the records of all arcs incoming in  $r_z$  as well as all records of the auxiliary file storing whether an arc encoding has been processed. Moreover during each iteration, such auxiliary file is written. Therefore the I/O complexity of an iteration of the while loop regarding the arc incoming in  $r_z$  is equal to 3 times the number of arcs of  $G_O$  incoming in  $r_z$ .

Now we have to determine the number of iterations of the while loop at lines 5–15. A consequence of Lemmas 28, 29, and 30 is that the condition at line 5 becomes false (and we exit from the while loop) only when all arcs in  $D_z$  are inserted in  $C$  in some iteration. The condition at line 10 that an arc encoding  $e$  is added to  $C$  only if  $|C| < M$  and  $e$  is not transitive. Therefore only the last iteration can terminate with a set  $C$  containing fewer than  $M$  elements. Hence the number of iterations is equal to  $\lceil |D_z|/M \rceil$ . Consequently the I/O complexity of an iteration of the for loop over all reads in  $R$  (lines 3–15) is equal to  $3\lceil |D_z|/M \rceil |E_O(r_z)|$ , where  $E_O(r_z)$  is the set of arcs of  $G_O$  that are incoming in  $r_z$ .

Summing over all iterations of the for loop at lines 3–15 immediately proves the theorem.  $\square$

## 5 Conclusions

The first contribution of this paper is a compact representation of the overlap graph and of the string graph via string-intervals. More precisely, we have shown how a string-interval can be used to represent the set of reads sharing a common prefix, with a possible reduction in the overall space used.

Then, we have proposed the first known external-memory algorithm to compute the overlap graph, showing that it reads at most  $(3 + 6l)n$  records, where  $n$  is the total length of the input and  $l$  is the maximum length of each input string, using only a constant amount of main memory. A fundamental technical contribution is the improvement of

the CompleteExtensions procedure that has been introduced in [13] to compute, with a single scan of the BWT, all backward  $\sigma$ -extensions of a set of disjoint string-intervals. Our improvement allows to extend a generic set of string-intervals.

Finally, we have described a new external-memory algorithm for reducing an overlap graph to obtain the corresponding string graph, reading or writing  $3|E(G_O)|\lceil d/M \rceil$  records, where  $E(G_O)$  is the arc set of  $G_O$  and  $d$  is the maximum indegree of  $G_S$ , while using an amount of main memory necessary to store  $M/5$  integers (as well as some constant-sized data structure).

There are some open problems that we believe are interesting. The analysis of the algorithm complexity is not very detailed. In fact, we conjecture that some clever organization of the records and a more careful analysis will show that the actual I/O complexity is better than the one we have shown in the paper.

Another direction is to assess the actual performance of the algorithm on data originating from a set of sequences, such as those coming from transcriptomics [8, 19] or metagenomics [23], especially to verify the gain in disk usage.

## Acknowledgements

The authors acknowledge the support of the MIUR PRIN 2010-2011 grant “Automi e Linguaggi Formali: Aspetti Matematici e Applicativi” code 2010LYA9RH, of the Cariplo Foundation grant 2013-0955 (Modulation of anti cancer immune response by regulatory non-coding RNAs), of the FA 2013 grant “Metodi algoritmici e modelli: aspetti teorici e applicazioni in bioinformatica” code 2013-ATE-0281, and of the FA 2014 grant “Algoritmi e modelli computazionali: aspetti teorici e applicazioni nelle scienze della vita” code 2014-ATE-0382.

## References

- [1] M. I. Abouelhoda, S. Kurtz, and E. Ohlebusch. Replacing suffix trees with enhanced suffix arrays. *J. of Discrete Algorithms*, 2(1):53–86, Mar. 2004.
- [2] A. Aggarwal and J. S. Vitter. The input/output complexity of sorting and related problems. *Commun. ACM*, 31(9):1116–1127, 1988.
- [3] F. Alizadeh, R. M. Karp, L. A. Newberg, and D. K. Weisser. Physical Mapping of Chromosome: A Combinatorial Problem in Molecular Biology. *Algorithmica*, 13:52–76, 1995.
- [4] F. Alizadeh, R. M. Karp, D. K. Weisser, and G. Zweig. Physical Mapping of Chromosomes Using Unique Probes. *Journal of Computational Biology*, 2:159–184, 1995.
- [5] A. Bankevich, S. Nurk, D. Antipov, et al. SPAdes: A new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.*, 19(5):455–477, 2012.



- [6] M. J. Bauer, A. J. Cox, and G. Rosone. Lightweight Algorithms for constructing and inverting the BWT of string collections. *Theoretical Computer Science*, 483:134–148, Apr. 2013.
- [7] D. A. Benson, K. Clark, I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and E. W. Sayers. Genbank. *Nucleic Acids Research*, 42(D1):D32–D37, 2014.
- [8] S. Beretta, P. Bonizzoni, G. Della Vedova, Y. Pirola, and R. Rizzi. Modeling alternative splicing variants from RNA-Seq data with isoform graphs. *J. Comput. Biol.*, 16(1):16–40, 2014.
- [9] A. Blum, T. Jiang, M. Li, J. Tromp, and M. Yannakakis. Linear approximation of shortest superstrings. *J. ACM*, 41:630–647, 1994.
- [10] M. Burrows, D. J. Wheeler, M. Burrows, and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical report, Digital Systems Research Center, 1994.
- [11] Y. Chen, G. Dong, J. Han, B. W. Wah, and J. Wang. Multi-dimensional regression analysis of time-series data streams. In *Proceedings of the 28th international conference on Very Large Data Bases*, pages 323–334. VLDB Endowment, 2002.
- [12] R. Chikhi and G. Rizk. Space-efficient and exact de bruijn graph representation based on a bloom filter. *Algorithms for Molecular Biology*, 8:22, 2013.
- [13] A. Cox, T. Jakobi, G. Rosone, and O. Schulz-Trieglaff. Comparing DNA sequence collections by direct comparison of compressed text indexes. In *Algorithms in Bioinformatics*, volume 7534 of *LNCIS*, pages 214–224. Springer, Berlin, Germany, 2012.
- [14] A. J. Cox, M. J. Bauer, T. Jakobi, and G. Rosone. Large-scale compression of genomic sequence databases with the BurrowsWheeler transform. *Bioinformatics*, 28(11):1415–1419, June 2012.
- [15] C. Demetrescu, I. Finocchi, and A. Ribichini. Trading off space for passes in graph streaming problems. *ACM Transactions on Algorithms*, 6:1, Jan. 2009.
- [16] R. Diestel. *Graph Theory*. Graduate Texts in Mathematics. Springer-Verlag, Heidelberg, third edition, 2005.
- [17] P. Ferragina and G. Manzini. Indexing compressed text. *J. ACM*, 52(4):552–581, 2005.
- [18] M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on data streams. In J. M. Abello and J. S. Vitter, editors, *External Memory Algorithms*, volume 50 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, chapter Computing on Data Streams, pages 107–118. American Mathematical Society, Boston, MA, USA, 1999.

- [19] V. Lacroix, M. Sammeth, R. Guigo, and A. Bergeron. Exact transcriptome reconstruction from short sequence reads. In K. Crandall and J. Lagergren, editors, *Algorithms in Bioinformatics*, volume 5251 of *Lecture Notes in Computer Science*, pages 50–63. Springer Berlin Heidelberg, 2008.
- [20] T. Lam, R. Li, A. Tam, S. Wong, E. Wu, and S. Yiu. High throughput short read alignment via bi-directional BWT. In *Bioinformatics and Biomedicine, 2009. BIBM '09. IEEE Int. Conf. on*, pages 31–36, 2009.
- [21] A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K. Garimella, D. Altshuler, S. Gabriel, M. Daly, and M. A. DePristo. The genome analysis toolkit: a mapreduce framework for analyzing next-generation dna sequencing data. *Genome research*, 20(9):12971303, September 2010.
- [22] E. Myers. The fragment assembly string graph. *Bioinformatics*, 21(suppl. 2):ii79–ii85, 2005.
- [23] J. Pell, A. Hintze, R. Canino-Koning, A. Howe, J. M. Tiedje, and C. T. Brown. Scaling metagenome sequence assembly with probabilistic de Bruijn graphs. *ArXiv e-prints*, June 2012.
- [24] Y. Peng, H. C. Leung, S.-M. Yiu, and F. Y. Chin. Idba-ud: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*, 28(11):1420–1428, 2012.
- [25] G. Rosone and M. Sciortino. The burrows-wheeler transform between data compression and combinatorics on words. In P. Bonizzoni, V. Brattka, and B. Lwe, editors, *The Nature of Computation. Logic, Algorithms, Applications*, volume 7921 of *Lecture Notes in Computer Science*, pages 353–364. Springer Berlin Heidelberg, 2013.
- [26] F. Shi. Suffix arrays for multiple strings: A method for on-line multiple string searches. In J. Jaffar and R. Yap, editors, *Concurrency and Parallelism, Programming, Networking, and Security*, volume 1179 of *Lecture Notes in Computer Science*, pages 11–22, Berlin, Germany, 1996. Springer Berlin Heidelberg.
- [27] J. Simpson and R. Durbin. Efficient construction of an assembly string graph using the FM-index. *Bioinformatics*, 26(12):i367–i373, 2010.
- [28] J. Simpson and R. Durbin. Efficient de novo assembly of large genomes using compressed data structures. *Genome Res.*, 22:549–556, 2012.
- [29] J. Simpson, K. Wong, S. Jackman, et al. ABySS: a parallel assembler for short read sequence data. *Genome Res.*, 19(6):1117–1123, 2009.
- [30] L. G. Valiant. General purpose parallel architectures. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science (Vol. A)*, pages 943–973. MIT Press, Cambridge, MA, USA, 1990.

- [31] J. Vitter and E. Shriver. Algorithms for parallel memory, i: Two-level memories. *Algorithmica*, 12(2-3):110–147, 1994.
- [32] J. S. Vitter. External memory algorithms and data structures: Dealing with massive data. *ACM Comput. Surv.*, 33(2):209–271, June 2001.